

**RULE-BASED APPROACH FOR DETECTING  
BOTNET BASED ON DOMAIN NAME SYSTEM**

**KAMAL IBRAHIM AHMED ALIEYAN**

**UNIVERSITI SAINS MALAYSIA**

**2018**

# **RULE-BASED APPROACH FOR DETECTING BOTNET BASED ON DOMAIN NAME SYSTEM**

**by**

**KAMAL IBRAHIM AHMED ALIEYAN**

**Thesis submitted in fulfilment of the requirements  
for the degree of  
Doctor of Philosophy**

**April 2018**

## DEDICATION

*To my appreciated father "Ibrahim Ahmed Alieyan"*

*To my dearest mother "Fatimah Ahmed Al Muhdi"*

*To my beloved wife "Tamara Mahmoud Gear"*

*To my lovely kids "Ibrahim, Rimas, Osied"*

*To my dearest Family "Mohamad, Ahmad, Hamza, Anas, Amal,  
Manal, Asma"*

## ACKNOWLEDGEMENT

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ  
{نَرْفَعُ دَرَجَاتٍ مِّنْ نَّشَأٍ وَفَوْقَ كُلِّ ذِي عِلْمٍ عَلِيمٌ - 76 - } , {سورة يوسف}

First of all, I would like to thank “Allah”, our creator, for giving me His blessings and the opportunity to complete my PhD.

Now, I would like to share my happiness by acknowledging some of the numerous people who have helped me directly or indirectly in shaping my academic career over the past years. I would like to express my deepest gratitude to my learned supervisors Dr. Mohammed F. Anbar (main supervisor), and Dr. Ammar ALmomani (field-supervisor), for their encouragement, innovative suggestions, and invaluable help during the entire period of my PhD at National Advanced IPv6 Centre of Excellence (NAv6), which is a high-profile organization. I also wish to thank my research committee members for providing insightful and constructive comments.

Most importantly, I would like to thank my wife Tamara. Her support, encouragement, patience, and unwavering love were undeniably the foundation upon which the past nine years of my life have been built. Her tolerance of my occasional discourteous moods is a testament to her unyielding devotion and love. I thank my parents, Ibrahim and Fatimah, for their faith in me and for allowing me to be as ambitious as I wished. It was under their watchful eye that I gained the drive and ability to tackle challenges head-on. I would also like to thank my brothers and sisters and Dr. Mohammad Aluthman for providing unconditional support and guidance as needed. Last but not the least, I wish to dedicate this work to my kids Ibrahim, Rimas, and Osaid.

## TABLE OF CONTENTS

Acknowledgement .....	ii
Table of Contents .....	iii
List of Tables .....	vii
List of Figures .....	ix
List of Abbreviations.....	xi
Abstrak .....	xii
Abstract .....	xiv
<b>CHAPTER 1 INTRODUCTION.....</b>	<b>1</b>
1.1 Background.....	2
1.1.1 Botnet .....	2
1.1.2 Domain Name System (DNS) .....	3
1.1.3 DNS-based Botnet Detection Approaches .....	4
1.2 Research Motivation.....	4
1.3 Research Problem .....	5
1.4 Research Objectives.....	7
1.5 Research Scope .....	7
1.6 Research Contributions.....	8
1.7 Research Steps .....	8
1.8 Thesis Organization .....	11
<b>CHAPTER 2 LITERATURE REVIEW.....</b>	<b>12</b>
2.1 Introduction.....	12

2.2	Background.....	12
2.2.1	Growth of Botnet.....	13
2.2.2	Overview of Domain Name System (DNS).....	14
2.2.3	Botnet Life Cycle .....	15
2.3	Related Work.....	18
2.3.1	HoneyNet-based Approaches.....	18
2.3.2	Intrusion Detection System (IDS).....	20
2.3.3	Summary and Discussion of Related Works.....	42
2.4	Feature Selection Approaches .....	45
2.4.1	Feature Ranking by Information Gain Ratio (IGR) Algorithm .....	46
2.4.2	Feature Ranking Principal Component Analysis (PCA) Algorithm .....	47
2.5	Summary.....	50

## **CHAPTER 3    RULE-BASED APPROACH FOR DETECTING BOTNET**

	<b>BASED ON DNS .....</b>	<b>51</b>
3.1	Introduction.....	51
3.2	Overview of the Proposed Approach (RADBDNS).....	51
3.2.1	Data Pre-processing.....	52
3.2.2	DNS Features Selection .....	52
3.2.3	DNS-based Botnet Detection .....	53
3.3	Requirements of the Proposed RADBDNS Approach .....	53
3.4	Proposed Approach (RADBDNS Approach) .....	53
3.4.1	Dataset Pre-Processing (Stage 1) .....	56
3.4.2	DNS Features Selection (Stage 2).....	60
3.4.3	DNS -Based Botnet Detection (Stage 3) .....	63

3.5	Summary.....	76
-----	--------------	----

## **CHAPTER 4    DESIGN AND IMPLEMENTATION OF THE PROPOSED**

### **RADBDNS APPROACH..... 77**

4.1	Introduction.....	77
4.2	Tools and Programming Languages for Implementation .....	77
4.2.1	Java Programming Language .....	77
4.2.2	Wireshark .....	77
4.2.3	MySQL Database .....	79
4.2.4	Weka.....	79
4.2.5	Experiment Environment .....	80
4.3	The Analysed DNS-based Botnets .....	82
4.3.1	Zeus .....	83
4.3.2	Citadel .....	83
4.3.3	Conficker.....	83
4.3.4	Waledac .....	83
4.4	Design of the Proposed Approach .....	84
4.4.1	Design of Data Collection and Pre-Processing Stage .....	84
4.4.2	Design of DNS Features Selection Stage.....	87
4.4.3	Design of DNS-based Botnet Detection Stage.....	90
4.4.4	Design of Rules-Based Behavior Detection.....	96
4.5	Summary.....	96

## **CHAPTER 5    EXPERIMENTAL RESULTS AND DISCUSSIONS ..... 98**

5.1	Introduction.....	98
-----	-------------------	----

5.2	Benchmark Datasets .....	98
5.2.1	ISOT Dataset (2010) .....	99
5.2.2	Network Information Management and Security (NIMS) Dataset .....	102
5.2.3	Hardware for The Proposed Approach.....	104
5.2.4	Software Specifications for The Proposed Approach .....	105
5.3	Evaluation Metrics.....	105
5.4	Ground Truth Test Scenarios.....	106
5.4.1	Ground Truth in ISOT Dataset (Scenario 1) .....	106
5.4.2	Ground Truth Test in NIMS Dataset (Scenario 2) .....	112
5.5	Comparison with Existing Approach.....	122
5.6	Summary.....	124

## **CHAPTER 6 CONCLUSION AND FUTURE WORK ..... 126**

6.1	Overview.....	126
6.2	Conclusion .....	126
6.3	Future Work.....	128

## **REFERENCES .....129**

## **APPENDICES**

## **LIST OF PUBLICATIONS**



## LIST OF TABLES

	<b>Page</b>
Table 1.1      Research Scope .....	8
Table 2.1      Summary of DNS Signature-based Botnet Detection Approaches .....	23
Table 2.2      DNS Host-based Botnet Detection Approaches .....	25
Table 2.3      Summary of Active Botnet Detection Approaches Based on DNS Traffic Analysis .....	27
Table 2.4      Summary of Botnet Detection Approaches Based on DNS Traffic .....	31
Table 2.5      Summary of Passive Statistical Approaches for Botnet Detection Based on DNS Traffic .....	33
Table 2.6      Summary of Passive Botnet Detection Approaches Based on DNS .....	41
Table 3.1      List of Basic Features for the Proposed Approach .....	58
Table 3.2      Alphabet Symbols and Their Frequencies .....	66
Table 3.3      Features Related to the Suspicious Domains .....	67
Table 3.4      Sample of Domains_log for DNS Queries .....	68
Table 3.5      Failure Connection Features .....	69
Table 3.6      Sample of Failed_DNS_Response_Log for DNS Responses .....	70
Table 4.1      Filters Used in Wireshark for Extracting the Basic Features .....	84
Table 4.2      New Features Extracted Using RADNSBD .....	90
Table 5.1      Summary of ISOT Dataset .....	100
Table 5.2      List of ISOT Malicious/Non-Malicious IP Address and Labeling Machines .....	101
Table 5.3      Total Number of DNS Packets .....	102
Table 5.4      NIMS Dataset Specifications .....	103
Table 5.5      Evaluation Metrics .....	105
Table 5.6      Samples of DNS Packets in Some Groups .....	107

Table 5.7	Sample of Number of Distinct IP Addresses in Each Group.....	107
Table 5.8	Source IP Addresses in Each Group .....	108
Table 5.9	Evaluation of RADBDNS Approach in ISOT Dataset.....	108
Table 5.10	IP Addresses and Accessed Domains in Group 230.....	110
Table 5.11	Abnormality DNS Response for Each Source IP In Group 230 And Group 233 .....	111
Table 5.12	Group Details for NIMS Dataset .....	113
Table 5.13	Distinct IPs in Alexa Trace.....	113
Table 5.14	Distinct IPs in Zeus Trace.....	113
Table 5.15	Distinct IPs in Citadel Trace.....	114
Table 5.16	Distinct IPs Addresses in Conficker Trace .....	114
Table 5.17	IP Address in the Selected Group Traces .....	115
Table 5.18	Evaluation of RADBDNS Approach Using NIMS Dataset .....	115
Table 5.19	Sample of Data in Alexa Trace (Non-Malicious) for Group 12 .....	117
Table 5.20	Sample of Data in for Group 303, Group 27, Group 2 .....	117
Table 5.21	Abnormality of DNS Response Messages for Source IP in Botnets Traces .....	120
Table 5.22	Abnormality of DNS Response Messages for Source IP in Alexa Trace.....	122
Table 5.23	RADBDNS vs PsyBoG in Term of Accuracy Detection and False Positive .....	123

## LIST OF FIGURES

	<b>Page</b>
Figure 1.1 Botnet Representation Issue in 2015, Source: (Arbor, 2015).....	2
Figure 1.2 Botnet Activity.....	3
Figure 1.3 Research Steps .....	10
Figure 2.1 Botnet Life Cycle.....	17
Figure 2.2 Classification of Botnet Detection Approaches Based on DNS Traffic Characteristics .....	19
Figure 2.3 Overview of the Malicious Fast-flux Service Networks Detection (Perdisci <i>et al.</i> , 2009) .....	35
Figure 2.4 Spatial Snapshot Fast Flux Detection System (SSFD) (Huang <i>et al.</i> , 2010) .....	37
Figure 2.5 Overview of EXPOSURE System (Bilge <i>et al.</i> , 2011) .....	38
Figure 2.6 Fuzzy Pattern Recognition Algorithm for Identifying Botnet Domain Names and IP Addresses (Wang <i>et al.</i> , 2014) .....	40
Figure 3.1 General Stages of Proposed Approach .....	52
Figure 3.2 Proposed Approach.....	54
Figure 3.3 Dataset Pre-Processing Stage .....	56
Figure 3.4 DNS Packet Structure (Liu and Albitz, 2006) .....	57
Figure 3.5 Result of Intersection.....	63
Figure 3.6 Flowchart of Checking the Domain Name .....	65
Figure 3.7 Cluster of Domains .....	72
Figure 4.1 Snapshot of Wireshark Tool .....	78
Figure 4.2 Snapshot of WEKA Tool.....	80
Figure 4.3 Experimental Steps of RADBDNS.....	81
Figure 4.4 Basic Features and Selected Features .....	87
Figure 4.5 Weka Snapshot of Fields Ranking Output Using IGR .....	88
Figure 4.6 Weka Snapshot of Fields Ranking Output Using PCA .....	89

Figure 4.7	Classification of DNS Traffic .....	91
Figure 4.8	Design of the Abnormality in DNS Queries .....	93
Figure 4.9	Design of the Abnormality in DNS Responses (First Step) .....	94
Figure 4.10:	Design of the Abnormality in DNS Responses (Second Step).....	95
Figure 5.1	Topology of ISOT Dataset.....	100
Figure 5.2	Flow Exporting Mechanism Generated by NIMS (Haddadi and Zincir-Heywood, 2016) .....	104
Figure 5.3	Average Entropy for Domain Names for Groups 230 and 233 ...	109
Figure 5.4	Average Entropy for Domain Names for each Groups in NIMS Dataset .....	116

## **LIST OF ABBREVIATIONS**

DNS	Domain Name System
DDoS	Distributed Denial of Service
C&C Server	Command and Control Server
IGR	Information Gain Ratio
PCA	Principal Component Analysis
SQL	Structured Query Language
QName	Query Name
QTYPE	Query Type
QCLASS	Query Class
TTL	Time To Live
DGA	Domain Generation Algorithm
NXDOMAIN	None Exist Domains
IRC	Internet Relay Chat
IP	Internet Protocol
GB	Giga Byte
IDS	Intrusion Detection System
DDNS	Dynamic Domain Name System
PSD	Power Spectral Density
SVM	Support Vector Machine
RDNS	Recursive Domain Name System
TCP	Transmission Control Protocol

# **PENDEKATAN BERASASKAN PERATURAN UNTUK MENGESAN BOTNET MENGGUNAKAN PELAYAN NAMA DOMAIN**

## **ABSTRAK**

Botnet merupakan suatu cabaran serius yang dihadapi oleh Internet sejak kebelakangan ini, membawa kerugian ekonomi kepada organisasi dan individu. Botnet terdiri daripada ribuan hos terjangkit yang menerima arahan daripada pelayan kawal dan perintah (C&C) yang dikendalikan oleh seseorang individu. Mengikut tradisi, pelayan saling bual Internet (IRC) digunakan sebagai pelayan C&C; mereka berkomunikasi dengan botnet melalui saluran IRC. Ini selalunya menyebabkan pentadbir rangkaian menyekat trafik IRC di rangkaian tersebut. Trend terbaru dalam pengeralahan tugas botnet melibatkan penggunaan saluran komunikasi alternatif, seperti pelayan sistem nama domain, antara pelayan C&C dan hos terjangkit (bot). Penggunaan saluran komunikasi alternatif telah membolehkan botnet untuk memintas tapisan rangkaian yang sedia ada. Tambahan, saluran-saluran ini tidak boleh disekat dengan mudahnya kerana trafik IRC adalah penting untuk aktiviti biasa rangkaian. Lebih lagi, botnet baru, seperti Conficker, Zeus dan Citadel telah menggunakan fluks-cepat DNS untuk mengelakkan pengesanan dan mengelakkan penyelidik daripada mencari dan menutup pelayan C&C. Oleh itu, tesis ini mencadangkan pendekatan berasaskan peraturan untuk mengesan botnet berasaskan DNS (RADBDNS), yang dapat meningkatkan kejituan pengesanan botnet yang berasaskan kepada trafik DNS. Pendekatan ini adalah berasaskan kepada peraturan yang bergantung kepada tingkahlaku pertanyaan dan maklum-balas DNS. RADBDNS terdiri daripada tiga peringkat: (1) pra-pemprosesan data untuk menapis trafik DNS daripada trafik rangkaian, dengan itu memberi fokus kepada sebahagian kecil trafik rangkaian tersebut untuk

mengurangkan overhead sistem ; (2) Pemilihan ciri-ciri DNS untuk memilih ciri-ciri yang paling bererti yang menyumbang kepada pengesanan botnet berasaskan DNS; dan (3) Pengesanan botnet berasaskan DNS untuk mengesan tingkah-laku abnormal pertanyaan dan maklum balas DNS dengan mengaplikasikan peraturan yang dicadangkan; hos yang memperlihatkan tingkah-laku abnormal pertanyaan dan maklum-balas DNS akan dikenalpasti sebagai bot. Pendekatan yang dicadangkan telah dinilai dengan menggunakan dua dataset tanda aras dalam dua senario. Dalam senario pertama, pendekatan yang dicadangkan diaplikasikan ke atas dataset ISOT, yang mengandungi trafik botnet serta trafik normal. Dalam senario kedua, pendekatan yang dicadangkan diaplikasikan ke atas dataset NIMS, yang mengandungi trafik botnet dan trafik normal secara berasingan. Keputusan penilaian menunjukkan pendekatan yang dicadangkan dapat mengesan botnet berasaskan DNS dengan kejituan sebanyak 99.66% dengan kadar positif palsu sebanyak 0.23%. Sementara itu, keberkesanan pendekatan yang dicadangkan itu dinilai dengan membandingkannya dengan pendekatan-pendekatan DNS yang lebih dikenali dan keputusan yang didapati menunjukkan bahawa pendekatan yang dicadangkan mempunyai prestasi yang lebih baik berbanding dengan pendekatan lain.

# **RULE-BASED APPROACH FOR DETECTING BOTNET BASED ON DOMAIN NAME SYSTEM**

## **ABSTRACT**

Botnets are a serious problem in today's Internet, and they result in economic damage for organizations and individuals. Botnets consist of thousands of infected hosts that receive instructions from command and control (C&C) servers operated by an individual. Traditionally, Internet Relay Chat (IRC) servers are used as C&C servers and communicate with the botnet through IRC channels. As a result, network administrators often block IRC traffic on their networks. Recent trends in botnet development have seen the use of alternative communication channels, such as domain name server (DNS), between the C&C servers and infected hosts (bots). The use of alternative communication channels has allowed botnets to bypass common network filters. Furthermore, these channels cannot be blocked as simply as IRC traffic because they are essential for normal network activity. Recent botnets such as Conficker, Zeus, and Citadel have used DNS fast flux to avoid detection and to reduce the ability of researchers to find and shut down the C&C servers. Therefore, this thesis proposes a rule-based approach for detecting botnet based on DNS (RADBDNS), which can enhance the accuracy of detecting botnets based on DNS traffic. RADBDNS uses a rule based on DNS query and response behaviors, and it consists of the following three stages: (1) data pre-processing to filter DNS traffic from the network traffic in the datasets, (2) DNS feature selection to select the most significant features that contribute to the detection of the botnet based on DNS, and (3) DNS-based botnet detection, which aims to detect abnormal behavior of DNS queries and responses by applying the proposed rules on DNS queries and responses. The host that exhibits



abnormality in DNS queries or DNS responses will be identified as a bot. The proposed approach is evaluated using two benchmark datasets in two scenarios. The first scenario applies the proposed approach on the ISOT dataset, which contains botnet traffic along with legitimate traffic. The second scenario applies the proposed approach on the NIMS dataset, which contains legitimate traffic and malicious traffic separately. The result shows that the proposed approach can detect the botnet-based DNS with 99.66% accuracy and a false positive rate of 0.23%. The effectiveness of the proposed approach is evaluated through a comparison with well-known DNS-based approach. Results show that the proposed approach outperforms other approach.

## CHAPTER ONE

### INTRODUCTION

The internet has become an element of critical infrastructure that plays a vital role in people's lives, encompassing communication, education, government services, banking, and commerce. Naturally, these segments are targeted throughout the Internet by malicious attackers. One of the most prominent emerging malicious systems is the botnet. Botnets have become a serious threat to the web (AsSadhan *et al.*, 2017; Silva *et al.*, 2013). Using botnets, attackers can gain control of large-scale distributed networks associated with the infected machines to achieve their goals (Bilge *et al.*, 2014).

As botnets proliferate, they have been causing significant losses to the global economy. For example, a cyber-attack on Christmas Eve against the website of a financial institution in regional California helped to distract bank officials from an online account takeover against one of its clients, netting cyber thieves more than \$900,000 (Harris *et al.*, 2014). Another example is the use of botnets in online banking transactions, resulting in estimated losses of over 100 million dollars (e Silva, 2017). Furthermore, advertisers have lost at least \$346 million owing to the actions of less than 15% of the botnet population (Chen *et al.*, 2017b). According to the Arbor report (2015), botnets triggered 76% of the distributed denial of services (DDoS) attacks, as shown in Figure 1.1 (Arbor, 2015). The other botnet activity represents 47% of attacks. Thus, botnets represent a major security concern.

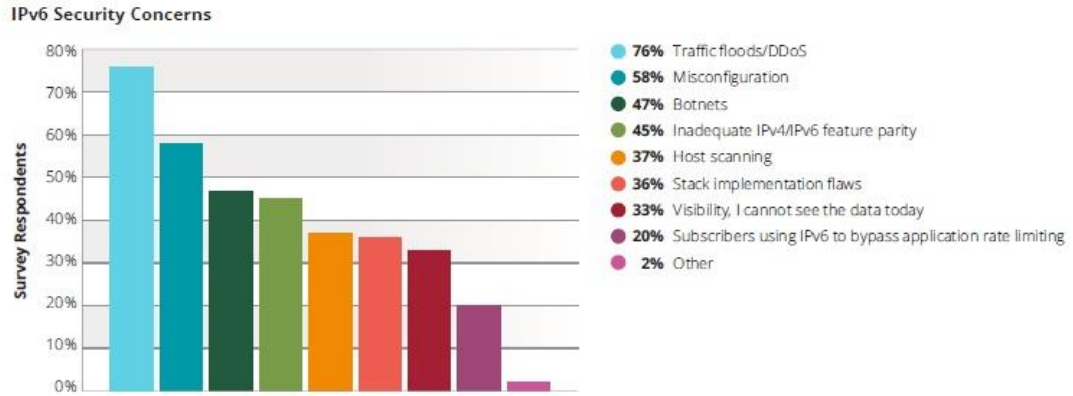


Figure 1.1: Botnet Representation Issue in 2015, Source: (Arbor, 2015)

## 1.1 Background

This Section provides an explanation of botnets, the domain name system (DNS), and DNS-based botnet detection techniques.

### 1.1.1 Botnet

A botnet is a group of infected computers on the Internet controlled by a botmaster through command and control (C&C) servers (Kwon *et al.*, 2016; Yu, 2014). Botnets perform many activities such as click fraud, spam, DDoS attacks, and hosting phishing pages, as shown in Figure 1.2. Botnets utilize a hidden communication channel to receive commands from their operator and communicate their current status (Tiirmaa-Klaar *et al.*, 2013b). Botnets sizes have increased dramatically in recent years, to the extent that networks of more than one million machines can launch cyber-attacks (Yukonhiatou *et al.*, 2014). These botnets need to control and manage all their distributed infected machines using a reliable system. One method to manage botnets is through the DNS.

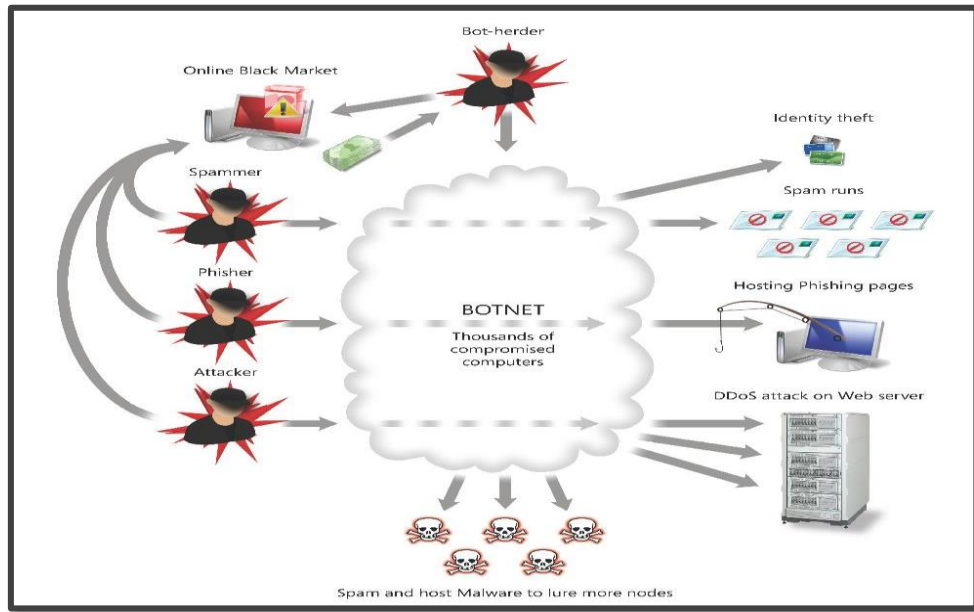


Figure 1.2: Botnet Activity

### 1.1.2 Domain Name System (DNS)

The DNS is a service provided by the Internet to translate domain names into internet protocol (IP) addresses (Agrawal *et al.*, 2014; Li *et al.*, 2017). The distributed and global nature of the DNS inspires cyber criminals to perpetrate attacks on a global scale. Consequently, the DNS is highly misused by attackers to manage their botnets (He *et al.*, 2010). Such attacks enable widespread exploitation by malicious systems and result in significant losses.

Recently, domain names have been used to operate malicious networks, such as botnets and other types of malicious software (malware). Studies have revealed that it is a challenge to keep track of malicious domains by web content analysis or human observation, because of the large number of domains (Davuth and Kim, 2013; Kwon *et al.*, 2016).

### 1.1.3 DNS-based Botnet Detection Approaches

There are two methods for discovering active botnets (Rahim and Bin Muhaya, 2010): active detection, which involves capturing live instances of running botnets, and passive detection, which involves capturing suspicious traffic from the network (Govil and Govil, 2007). Although many DNS-based botnet detection approaches have been proposed, such as (Antonakakis *et al.*, 2010; Antonakakis *et al.*, 2011), (Holz *et al.*, 2008), (Perdisci *et al.*, 2009), and (Bilge *et al.*, 2011; Bilge *et al.*, 2014), these approaches suffer from limitations, such as high false positive detection rates, and require the selection of many features. Additionally, the error in the classification process increases over time, especially when dealing with unknown botnet attack. Therefore, new approaches are needed to increase the detection of unknown botnets, because several proposed methods are unable to successfully detect these threats (García *et al.*, 2014).

## 1.2 Research Motivation

Researchers consider botnets to be critical tools for criminals, because of the low risk and high potential profit associated with their use. Botnet programs are growing faster than the response of security firms, and they represent a serious threat to the existence of the commercial Internet (e Silva, 2017).

Infections on the Internet caused by malware are estimated to account for 16–25% of global Internet traffic, which comes from communication between various types of malware (Camelo *et al.*, 2014). Botnets are used by criminals for gaining economic profit, as well as for politically motivated activities. Although many efforts have been made recently to mitigate the effectiveness of botnets, newer and more sophisticated versions are likely to re-emerge (Tiirmaa-Klaar *et al.*, 2013a).

The existing detection approaches of DNS-based botnets that use simple heuristics are inefficient, because modern botnets are sophisticated and encrypted. This has motivated researchers to suggest many approaches that can monitor DNS traffic to detect botnets. Such approaches include the following:

- Increasing the attention of DNS-based botnet detection approaches as efficient techniques in the Internet security and monitoring fields. These approaches can be implemented to detect the abnormal DNS traffic, which leads to the detection of DNS-based botnets.
- Needs to differentiate between DNS with trusted reputations and the domain names generated by botnets.

### **1.3 Research Problem**

The botnet phenomenon is one of the most significant threats to cyber security, providing cybercriminals with a distributed platform for several illegal activities, such as launching DDoS attacks, click fraud, phishing, identity theft, and spamming (Garg and Sharma, 2017; Kirubavathi and Anitha, 2014). According to McAfee Labs, the number of newly discovered malware samples reached approximately 11 million in the fourth quarter of 2012 (Shi *et al.*, 2013), and botnets represent a more dangerous threat to the internet (Chen and Lin, 2015). This illustrates the importance of developing efficient and useful detection approaches to control botnets.

DNS is a fundamental element of the functionality of the Internet, and the security of DNS affects the entire Internet (Shan *et al.*, 2014). Botnets use DNS as alternative communication channels between the C&C servers and infected hosts (bots) (Frosch *et al.*). The use of alternative communication channels has allowed botnets to bypass common network filters. Moreover, these channels cannot be

blocked as simply as IRC traffic, which has been used in traditional botnets, as they are essential for normal network activity. Furthermore, recent botnets such as Conficker have used DNS fast-flux to avoid detection and to reduce the ability of researchers to find and shut down the C&C servers (Davuth and Kim, 2013).

There are many approaches designed for botnet detection. The common detection approaches include signature-based detection, anomaly-based detection, mining-based detection, and DNS-based detection (Feily *et al.*, 2009; Rahim *et al.*, 2014). The DNS-based detection methods do not need specific knowledge about the botnet protocol and structure, and these can thus be a more efficient way of detecting botnets (Kang, 2011). Nevertheless, botnets are becoming more sophisticated, and resistant to detection. There are many approaches that can be used to evade detection, such as encrypted communication, domain generation algorithm (DGA), fast-flux, and double fast-flux (Camelo *et al.*, 2014).

The appropriate selection of features has critical positive impact on the performance of detection approaches. As a result, botnet detection approaches require optimization of the chosen feature set to capture the relevant botnet traffic heuristics (Stevanovic and Pedersen, 2013). Additionally, the existing DNS-based botnet detection approaches can detect real-world botnets with a considerably high false positive rate, which leads to a reduction in the accuracy of DNS-based botnet detection, because they do not take into consideration the features that significantly contribute to detecting botnets based on DNS. One of the main limitations of the existing approaches is the high false positive rate (Ji *et al.*, 2014).

The research problem can be summarized as follows:

- The existing approaches for detecting botnets based on DNS traffic do not consider the significant features in the DNS packet that can contribute to detecting botnets accurately.
- Most of existing approaches for detecting botnets based on DNS traffic suffer from high rates of false positive detection.

#### **1.4 Research Objectives**

The main goal of this research is to propose a rule-based approach for DNS-based botnet detection with improved accuracy. The following objectives are set to achieve the main goal of this research:

- To propose an ensemble feature selection mechanism to select the distinguished features to characterize the behavior of botnets based on DNS.
- To propose a set of new features that contribute to the detection of botnets based on DNS queries and DNS responses.
- To propose a rule-based mechanism to detect the DNS-based botnets.
- To evaluate the proposed approach and compare it with existing DNS-based botnet detection approach in terms of accuracy detection and false positive rate.

#### **1.5 Research Scope**

This research aims to differentiate between botnet-based DNS traffic and legitimate DNS traffic by focusing on the DNS query and response behaviors. The proposed approach is limited to detecting botnets based on DNS traffic in IPv4 and IPv6 environments with two benchmark datasets, as presented in Table 1.1.



Table 1.1: Research Scope

No	Items	Scope of Research
1.	Environment	IPv4 & Ipv6
2.	Attack type	Botnet based on DNS
3.	Target layer	Application layer (DNS)
4.	Features Selection	DNS Queries and DNS Responses
5.	Dataset	Benchmark dataset (ISOT, NIMS)
6.	DNS query and DNS response features	DNS query and DNS response features will be determined by analyzing Waledac, Zeus, Conficker, and Citadel DNS traffic
7.	Performance Metrics	Accuracy, False Positive

## 1.6 Research Contributions

The main contribution of this research is proposing a rule-based approach, which is designed to more accurately detect botnets based on DNS, called the rule-based approach for DNS-based botnet detection (RADBDNS). The contributions of this research are summarized as follows:

- An ensemble feature selection mechanism to select the distinguished features that can characterize botnets on based DNS.
- A set of new features that contribute the detection of botnets based on DNS.
- A rule-based mechanism to detect the DNS-based Botnet.

## 1.7 Research Steps

This research is conducted using a combination of theoretical analysis and experiments to examine the performance of security approaches that detect botnets based on DNS traffic using DNS query and response behaviors. Figure 1.3 depicts the complete research steps of this research. The proposed approach includes the detection of botnets based on DNS traffic. Thus, many botnets defeated the detection, and the

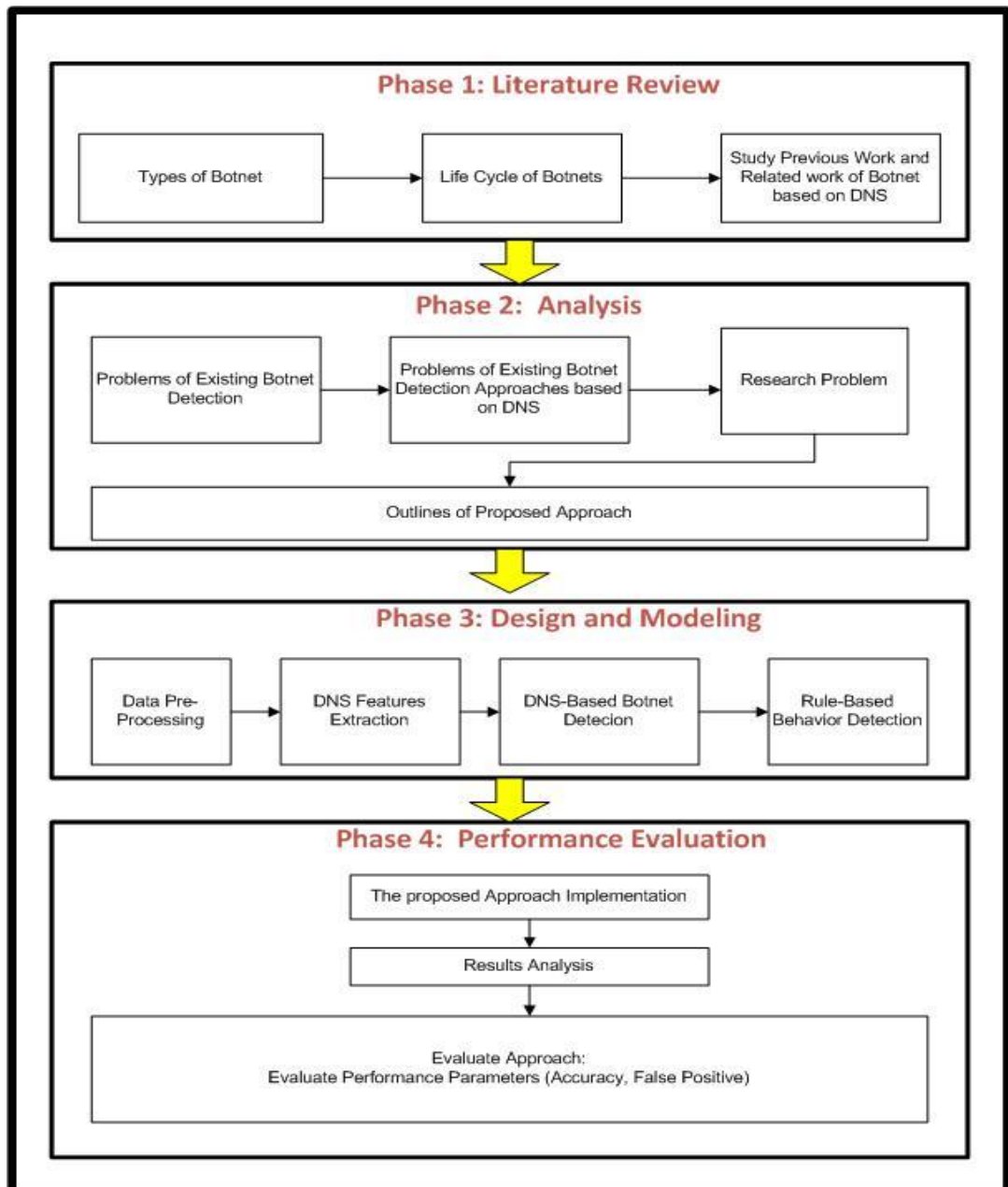
failure caused by the attacker can be minimized. As illustrated in Figure 1.3, there are four research steps to achieve our goals:

The **first phase** is the literature review. This step starts with an investigation of the types of botnet attack based on DNS, followed by a study of the life cycle of botnets, and a review of previous work and related botnet detection approaches based on DNS traffic.

The **second phase** is the analysis. This phase starts with an investigation of the problems of the existing botnet detection approaches. Then, the approaches used in botnet detection based on DNS are considered to identify the limitations of these approaches. As a result, the problem statement of this work is identified to outline the proposed approach.

The **third phase** is the design and modeling. This phase presents how to design the proposed approach and discusses the feature extraction for the proposed approach to enhance the detection performance, which works on two main behaviors of botnet based on DNS: abnormal DNS query behaviors and abnormal DNS response behaviors. Then the rule-based approach for behavior detection is discussed and presented.

The **fourth phase** is the performance evaluation. This phase discusses the implementation and experimental environment to evaluate the performance of the proposed approach. This phase starts with the approach implementation and then results in an analysis considering performance parameters such as the accuracy and false positive detection rate.



**Figure 1.3: Research Steps**

## 1.8 Thesis Organization

This thesis is structured into the following six chapters:

**Chapter 1** presents the background, research problem, research motivation, scope, objectives, and contributions of this research. This chapter also discusses the need for detection approaches for botnet based on DNS traffic.

**Chapter 2** discusses the background of the research and related studies. This chapter critically reviews the existing approaches for the detection of botnets based on DNS traffic and presents their advantages and limitations. Moreover, this chapter shows a new taxonomy for approaches of botnets based on DNS. Finally, this chapter comprehensively discusses the gaps in the existing approaches.

**Chapter 3** explains the methodology phases of the proposed approach for the detection of botnet based on DNS traffic within network traffic. Additionally, it describes the integrated phases of the proposed approach.

**Chapter 4** presents the design and tools used for the proposed approach. This chapter contains the feature selection and rules designs of the proposed approach. This chapter also explains the implementation of the phases in detail.

**Chapter 5** reports the experiments and their results. It also presents a comprehensive analysis of the results achieved using the proposed approach. In addition, this chapter evaluates the performance of the proposed approach in comparison with existing approaches.

**Chapter 6** presents the conclusions drawn from our work and suggests possible directions for future research.

## **CHAPTER TWO**

### **LITERATURE REVIEW**

#### **2.1 Introduction**

The previous chapter has discussed the importance of botnet detection approaches based on DNS. This chapter presents a detailed and comprehensive background of botnet detection, particularly botnet detection based on DNS traffic. It also presents an overview of most approaches used in detection of botnets based on DNS to provide an understanding of the proposed research in the field of botnet detection based on DNS traffic characteristics. In addition, this chapter provides a review of the most well-known DNS-based botnet detection approaches. The chapter is organized as follows. Section 2.2 provides a background of botnets, Section 2.3 discusses the related works in the botnet detection based on DNS, Section 2.4 formulates the feature selection for botnet detection based on DNS, and the chapter is summarized in Section 2.5.

#### **2.2 Background**

The increasing reliance on the Internet nowadays presents many challenges in terms of managing the Internet and the application usage, such as protecting user data, privacy, integrity, and availability. In the last few years, the Internet has played a central role in the world, especially in the sectors of communication, education, government services, and banking (Stevanovic *et al.*, 2012). Unfortunately, the increasing demand for applications become a threat to user privacy and data security (Stevanovic and Pedersen, 2013). This section discusses in detail the growth of botnets, an overview of DNS, and the botnet life cycle with botnets depending on DNS.

### 2.2.1 Growth of Botnet

A botnet is a software program that infects computers, known as bots, and manipulates them for malicious purposes. Bots run small scripts built to carry out specific automated tasks (Alomari *et al.*, 2012). These bots are controlled by one or a small collaborative group of attackers known as “botmasters” (Lu *et al.*, 2011b). Based on statistics from McAfee Labs (McAfee., 2015), the number of newly discovered malware samples reached 50 million in the fourth quarter of 2014 and is expected to reach 500 million by the end of 2015. Moreover, the Internet traffic consists of up to 80% of botnet traffic related to SPAM emails originating from known botnets such as Grum, Cutwail, and Rustock (Karim *et al.*, 2014). Currently, a large-scale botnet may consist of more than one million PCs launching cyber-attacks (Yukonhiatou *et al.*, 2014).

Botnets differ from other types of malware by utilizing communication channels to receive commands and report their current status to their operator(s) (Tiirmaa-Klaar *et al.*, 2013b). According to a 2013 FBI report, 10 international hackers were arrested for using botnets to steal more than \$850 million through a group of compromised computers, using stolen personal financial information (Hsu *et al.*, 2017).

In fact, botnets have specific characteristics as compared to other types of malware. For instance, a botmaster can control the infected machines and send commands without directly communicating with them. There are also many bots working in a coordinated way and taking instructions from the botmaster to instantiate coordinated attacks such as DDoS, spam distribution, and click fraud (Emre, 2011).

Additionally, botnets provide these frauds as services from the botnet operators, which are considered part of the botnet economy (Tiirmaa-Klaar *et al.*, 2013a).

Botnets are one of the most significant threats to cyber security, as they are considered a launching pad for several illegal activities such as DDoS, click fraud (Almomani *et al.*, 2013), phishing, identity theft (Al-Momani *et al.*, 2011), spamming (Kirubavathi and Anitha, 2014), phishing, and malware distribution (Zeidanloo *et al.*, 2010). Until now, there exists no permanent solution for the detection or mitigation of botnet threats, because their approaches and methods keep changing over time (Karim *et al.*, 2014).

### **2.2.2 Overview of Domain Name System (DNS)**

DNS is a fundamental element of the functionality of the Internet, which converts domain names to their corresponding IP addresses. However, the security of the DNS is the responsibility of the whole Internet collaboratively (Shan *et al.*, 2014). The distributed and global system of the DNS motivates the cyber criminals to attack on a global scale (Davuth and Kim, 2013). To commit their crimes, attackers make use of DNS services to operate malicious networks, such as botnets and other types of malware (He *et al.*, 2010).

In addition, studies have demonstrated the challenges in tracking malicious domains using web content analysis or human observation owing to the huge number of available domains within cyberspace (Davuth and Kim, 2013). Unfortunately, botnets use the DNS traffic as any other legitimate host, which makes differentiating legitimate DNS traffic from illegitimate traffic a very challenging problem (Manasrah *et al.*, 2009). Moreover, botnet owners attempt to hide their communication with the bots to obstruct any deployed botnet detection processes (Rodríguez-Gómez *et al.*,

2013). The attackers or botmasters use the DNS services to hide their C&C IP address to make the botnets reliable and easy to migrate from one server to another without being noticed (Choi *et al.*, 2007).

The use of DNS as an alternative communication channel has allowed botnets to bypass common network filters. Moreover, these channels cannot be blocked as easily as IRC traffic can be, as they are essential for regular network activity. Furthermore, recent botnets such as Conficker, Zeus, and Citadel have used DNS fast-flux to avoid detection and to reduce the ability of researchers to find and shut down the C&C servers (Haddadi *et al.*, 2015; Perdisci *et al.*, 2012).

### **2.2.3 Botnet Life Cycle**

The botnet detection process stands as an ongoing challenge for researchers and organizations. Therefore, understanding the botnet life cycle and their architecture may yield better detection mechanisms.

Generally, botnets apply a similar set of steps to recruit members and form the zombie army. These steps can be considered as the botnet life cycle. Figure 2.1 illustrates the steps of any botnet life cycle. The bots can typically be created and preserved in four phases, as shown in Figure 2.1.

#### **2.2.3(a) Exploitation Phase**

This phase is the first step in the botnet life cycle. The botmaster makes a remote infection by exploiting an existing vulnerability of software running on the victim host. The botmaster defrauds the victim user to execute malicious code on his machine, such as opening an email attachment (Abu Rajab *et al.*, 2006). In this phase, the bots need



to connect to a remote server to download the bot binaries. The connection to a remote server is established only after a DNS lookup command is issued by the compromised machine to map a domain name to its corresponding IP address (Abdullah *et al.*, 2013). This behavior of issuing a DNS lookup query is the dominant behavior of almost all botnets that exist in cyberspace (Manasrah *et al.*, 2009).

### **2.2.3(b) Rallying Phase**

In this phase, bots connect back to their botmaster through porting to a C&C server. The botmaster intends to make his botnet portable and stealth at the same time. Therefore, the botmaster equips his bots with a DNS lookup functionality to be able to perform DNS queries to locate the C&C server. Unfortunately, botmasters have learnt that a static IP address of the C&C is not effective, and vulnerable to identification and blacklisting. Therefore, they start to misuse the DNS services to hide the location of the C&C server behind a domain name rather than a static IP address. As a result, bots will rally to connect back to the C&C server as soon as they obtain the location of the required server (Choi *et al.*, 2007). This phase is considered vital to the success maintenance of the stealth nature and power of the botnet (Liu *et al.*, 2008).

### **2.2.3(c) Attack Execution Phase**

In this phase, the group of bots performs malicious activities on target machines as instructed by the botmaster through sending the required commands to the C&C servers. Bots will then fetch the command from the C&C server and start the malicious activity. For instance, the group of bots may receive commands from the C&C server to redirect users' requests to certain malicious websites by capturing the users' DNS queries (Rodríguez-Gómez *et al.*, 2013).

### 2.2.3(d) Update and Maintenance Phase

The last phase of the botnet life cycle is updating and maintaining the bots of the botnets. The botmaster needs to keep his bots up-to-date through instructing the bots to update their binaries from time to time for better coordination and patching (Karim *et al.*, 2014). Moreover, botmasters may require migrating his C&C server location frequently to evade the various detection measures (Karim *et al.*, 2014). Understanding this phase is very important, because botnets can be identified by observing the same network behaviors and communication patterns/frequency from the bots to their C&C server.

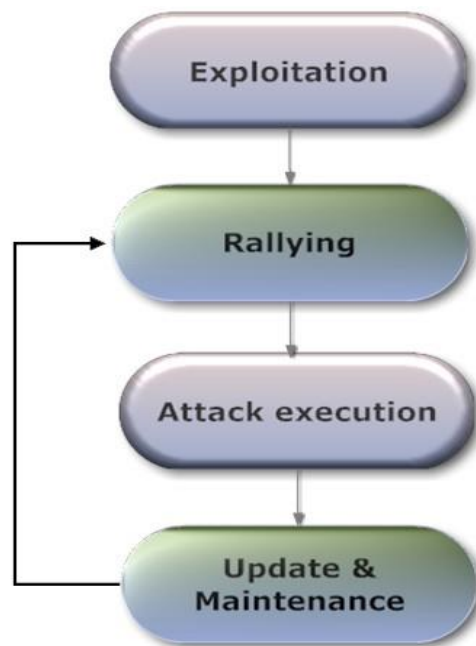


Figure 2.1: Botnet Life Cycle

## 2.3 Related Work

For the purpose of detecting botnets, different approaches were proposed to eliminate the danger of botnets attack. Moreover, researchers have produced different classifications for better understanding the botnet phenomenon and its structure (Feily *et al.*, 2009; Jing *et al.*, 2009; Karim *et al.*, 2014; Khattak *et al.*, 2013; Rodríguez-Gómez *et al.*, 2013; Silva *et al.*, 2013). The botnet detection approaches were mainly classified into two types: those based on installing and configuring a HoneyNet within the monitored network, and the IDS (Abdullah *et al.*, 2013; Jing *et al.*, 2009; Silva *et al.*, 2013; Zeidanloo *et al.*, 2010). A classification of botnet detection approaches based on analyzing DNS traffic is shown in Figure 2.2.

### 2.3.1 HoneyNet-based Approaches

HoneyNet-based approaches are used mainly to analyze and understand the behaviors and the characteristics of botnets. HoneyNets emulate known software and network vulnerabilities to be infected by botnets (Khattak *et al.*, 2013). HoneyNets are prepared to be self-contained and thwart the extension of botnets. In addition, HoneyNets are used to discover the capabilities of unknown attacks, the C&C system, the attackers' tools, approaches, and motivation (Karim *et al.*, 2014). Different approaches were proposed based on HoneyNet systems to capture botnet characteristics, such as in (Abu Rajab *et al.*, 2006; Dagon *et al.*, 2006; Li *et al.*, 2009b; Oberheide *et al.*, 2007; Pham and Dacier, 2011; Rieck *et al.*, 2010).

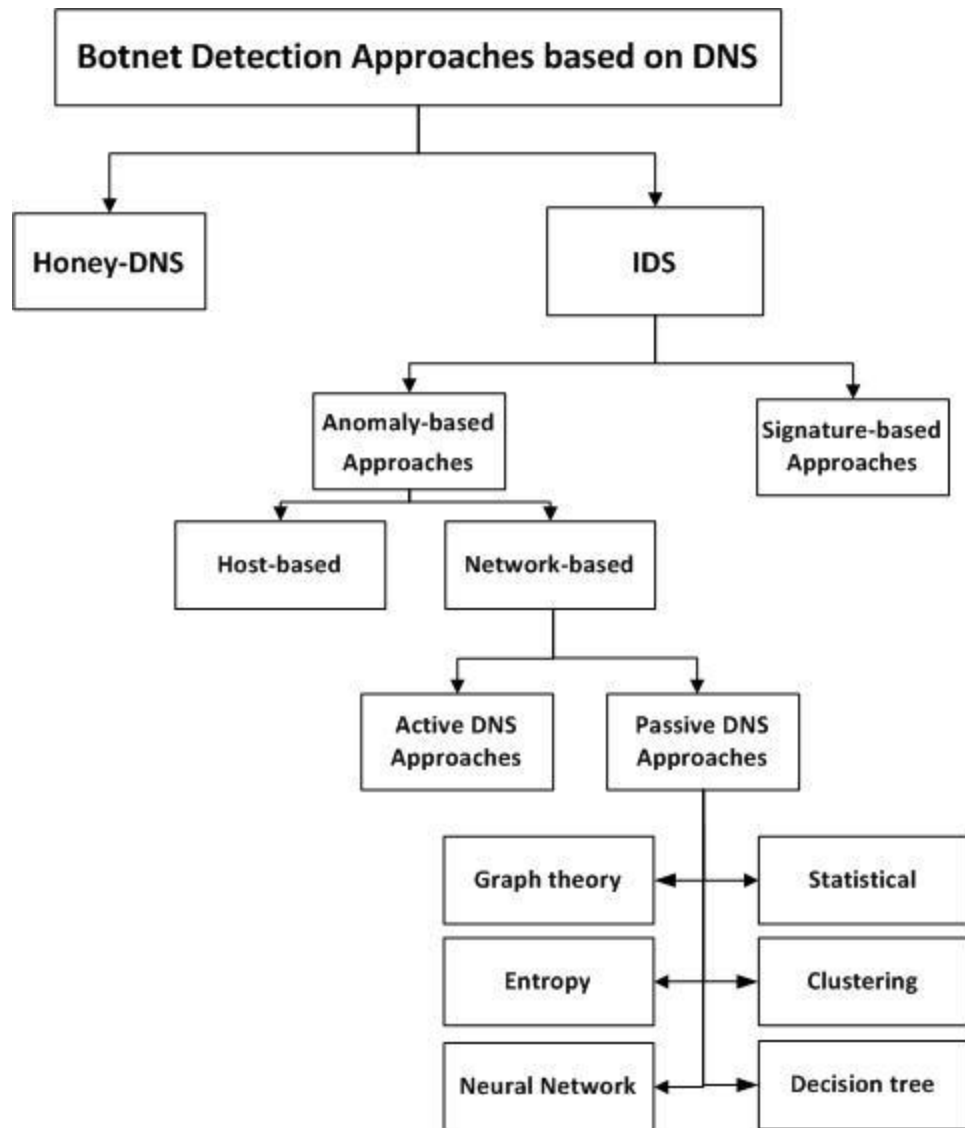


Figure 2.2: Classification of Botnet Detection Approaches Based on DNS Traffic Characteristics

HoneyNets are important to understanding botnet characteristics and technology (Silva *et al.*, 2013; Zeidanloo *et al.*, 2010). One of the known works on a HoneyNet that uses DNS queries is that proposed by (Oberheide *et al.*, 2007) who applies some basic statistics to the collected DNS queries. This work dealt with DNS queries targeting unused (i.e., darknet) address spaces and developed the concept of the honeydns system to assist honeypots to prevent the attackers from initiating their attacks (Aiello *et al.*, 2014a; Aiello *et al.*, 2014b).

HoneyNet-based systems are easy to build and deploy with minimum cost and resource requirements. Nevertheless, there are some drawbacks of the HoneyNet systems, including limited scalability and interaction with malicious activities. Moreover, attackers may use the honeyNet to learn new evasion approaches (Karim *et al.*, 2014). As a result, HoneyNets are aimed to recognize the features and mechanisms of botnets, but cannot detect bot infections all the time (Zeidanloo *et al.*, 2010).

### **2.3.2 Intrusion Detection System (IDS)**

The IDSs for botnet detection can be classified into two approaches: signature-based IDS (Abdullah *et al.*, 2013; Panimalar and Rameshkumar, 2014) and anomaly/behavior-based IDS (Gu *et al.*, 2007; Jing *et al.*, 2009; Li *et al.*, 2009a; Silva *et al.*, 2013; Vania *et al.*, 2013).

#### **2.3.2(a) Signature-based IDS**

Signature-based approaches detect only known bots through signature matching using an IDS detection system such as SNORT (SNORT). A DNS-based Black List (DNSBL) approach proposed by Ramachandran *et al.* (2006) is an example of a signature-based Botnet detection system (Ramachandran *et al.*, 2006). DNSBL-based approaches look for known bot signatures within the monitored DNS traffic. DNSBL-based approaches are also used to publish malicious and spamming activities online through collecting IP addresses of server machines or networks related to these activities. DNSBL based approaches attempt to recognize the botmasters' address and identify their location. However, the limitation of DNSBL-based approaches resides in maintaining an up-to-date database of known malicious addresses. Unfortunately,

one of the basic lines of defense against DNS abuses is domain name blacklisting (Oro *et al.*, 2010; Sinha *et al.*, 2008).

Similarly, Antonakakis *et al.* (2010) built a dynamic DNS (DDNS) reputation system called “Notos” that uses the passive DNS query data and analyzes the network and zone feature of a domain name. The Notos system assumed that malicious DNS queries have distinctive characteristics that are distinguishable from benign DNS queries (Antonakakis *et al.*, 2010). Thus, observing DNS queries and building models of known malicious and benign domains is feasible and might lead to a good result. A reputation score for the new domain observed was computed by models that give low scores for malicious domains and high scores for benign domains to differentiate between them. The Notos system has achieved high accuracy and low false positive rate, and it can recognize new domains before they get released to the public blacklist. However, the system needs a lot of history for a given domain name to reach a correct reputation score and it is inaccurate against frequently changing C&C domains, such as the hybrid botnet architecture that uses many master C&C nodes to distribute its commands (Kheir *et al.*, 2014).

In contrast to previous work, the mentor approach proposed by Kheir *et al.* (2014) depended on removing the legitimate domains from the botnet C&C blacklisted domains to reduce the false positive ratio during the detection process. The mentor approach implements scalable, positive DNS reputation approach that automatically removes benign or harmless domains recorded inside a blacklist of botnet C&C domains (Kheir *et al.*, 2014). The mentor approach collects statistical features about a suspected domain name, such as DNS properties and web content to build a DNS pruning model by applying a supervised learning approach into a labeled set of known benign and malicious domain names. The result of the mentor approach was effective

after testing over a public blacklist, and it removed benign domain names with a very low false positive rate. However, this approach relying on supervised machine learning and it also needs to be trained thus suffering from the ground truth problem (Stevanovic *et al.*, 2015).

Yadav *et al.* (2010) proposed an approach to detect “domain fluxes” in DNS traffic through searching for algorithmically generated patterns in domain names that are different from the domains generated by humans (Yadav *et al.*, 2010). They observed a distribution of alphanumeric characters together with the bigrams of all the domains mapped to the same set of IP addresses. However, this system is limited to the detection of C&C domains for only known malware samples that are correctly performed in the training phase, so it cannot identify unknown botnets (Kheir *et al.*, 2014). Table 2.1 shows a summary of DNS signature-based botnet detection approaches. Generally, the signature-based detection approaches have many limitations; mainly, they require constant updating to detect the new botnets or zero-day botnet attacks in which the signatures are evolving (Silva *et al.*, 2013).

Table 2.1: Summary of DNS Signature-based Botnet Detection Approaches

Proposed model	Mechanism	Weakness
DNSBL (Ramachandran <i>et al.</i> , 2006)	– Collecting published IP addresses of server machines	– Update DNS-based blacklist – Hard to design evasion approaches
Notos (Antonakakis <i>et al.</i> , 2010)	– Dynamic DNS reputation system – Uses passive DNS query data to analyze the network feature of a domain name	– Requires a long history for the given domain name to make reputation score – Unreliable with hybrid botnet
Mentor (Kheir <i>et al.</i> , 2014)	– Removing the legitimate domains from the botnet C&C domains blacklist	– Relying on supervised machine learning and it also needs to be trained thus suffering from the ground truth problem
Yadav (Yadav <i>et al.</i> , 2010)	– Detect domain fluxes in DNS traffic – Seeking for algorithmically generated patterns inherent to domain names	– Limited to known botnets – Attacks evaded detection during analysis

### 2.3.2(b) Anomaly DNS-based Botnet Detection Approaches

The anomaly- or behavior-based approaches attempted to detect botnets through analyzing network traffic for anomalies such as a sudden vast amount of traffic, traffic to unusual ports, high network latency, and anomalous behavior that may indicate the existence of bots in the network (Zeidanloo *et al.*, 2010). These approaches have the ability to identify new bots. The anomaly-based approach can be



categorized into host-based and network-based detection approaches (Karim *et al.*, 2014; Silva *et al.*, 2013; Zeidanloo *et al.*, 2010).

### **2.3.2(c) Host-based Anomaly Detection Approaches**

In host-based approaches, the monitoring and analyzing process are performed locally at each individual computer to detect any malicious activities through monitoring system processes, access to kernel level routines and system calls (Karim *et al.*, 2014). An example of the host-based approach is BotSwat, which was proposed by Stinson and Mitchell (2007) (Stinson and Mitchell, 2007). BotSwat focuses on the way bots respond to data received over the network through monitoring the execution of an arbitrary Win32 binary. However, the main limitation of the host-based botnet detection approach is that it is not scalable and limited to only bots within the monitored hosts. Moreover, to cover a wider view of the network, each individual host should be equipped with powerful monitoring tools that work collaboratively with others (Zeidanloo *et al.*, 2010).

One of the key approaches that focus on monitoring DNS traffic at the host/network level is the EFFORT framework that was proposed by (Shin *et al.*, 2012). This framework aims for effective and efficient detection through applying the multi module approach that correlates bot-related indications from different clients and network-level aspects. The EFFORT framework used a supervised machine learning algorithm to distinguish the queried domain names as being benign or malicious. The EFFORT framework is independent of topology, deployed communication protocol, and capable of detecting encrypted protocols. However, EFFORT is limited in its scope to botnets that depend on DNS services for identifying the address of their C&C